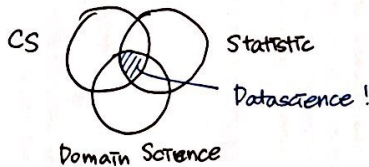


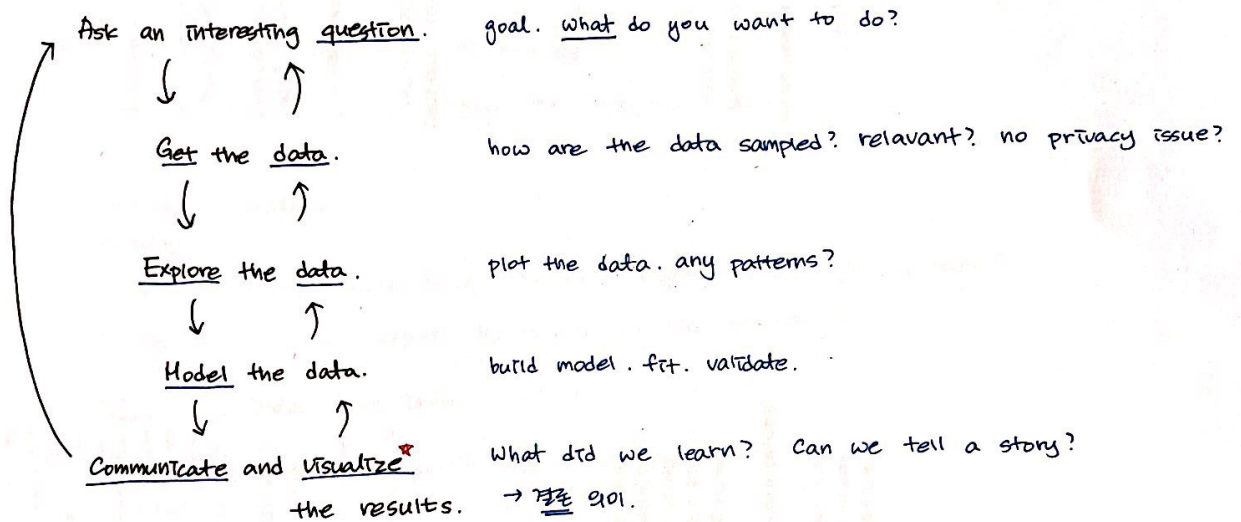
# Lecture 1. What is data science?

## Data science :

- Exploratory Data Analysis & Visualization
- ML & Statistics
- High-Performance Computing technologies for dealing with scale.



## Data science Pipeline



## Computer Science vs. Real Science

- |  |   |
|--|---|
| <ul style="list-style-type: none"> <li>- <u>algorithm</u> driven<br/>(focus on <u>methods</u>)</li> <li>- try to build their own clean virtual world<br/>(try to invent things)</li> <li>- random data can prove correctness of their methods</li> </ul> | <ul style="list-style-type: none"> <li>- <u>data</u> driven<br/>(focus on <u>results, findings</u>)</li> <li>- try to understand messy natural world<br/>(try to discover things)</li> <li>- "data" is important</li> </ul> |
|--|---|
- 
- 8/13 = 0.615

: care about what the number is!

8/13 ≈ 0.62

: care about what it means
- 
- |   |  |
|---|--|
| <ul style="list-style-type: none"> <li>- correct / wrong (정답 but)</li> <li>- <u>accuracy</u></li> </ul> | <ul style="list-style-type: none"> <li>- Nothing is completely wrong.</li> <li>- <u>meaning</u></li> </ul> |
|---|--|

## • Data Scientist

- : real scientist처럼 생각하는 필요가 없다.
- : are hired to produce insights ← come from understanding "meaning"

## • Developing Curiosity

- Data scientists should develop curiosity about the domain / application they are working in.
- broader perspective

## • Asking Good Questions

- 주어진 data set 으로부터 어떤 exciting 한 것을 할 수 있는가?
- 무엇을 알고 싶어하는가? (사람들, 나...)
- 어떤 테이아웃이 필요한가?

## • Practice Asking Questions : Who, What, Where, When, Why

### 1) Baseball-Reference.com : play의 통계 정보.

- Individual player's skill을 측정하는 가장 좋은 방법은 무엇?
- etc..

### 2) Demographic Questions

(인구통계학적인)

- Do left-handed people have shorter lifespans than right handers?
- Are heights and weights increasing in the population?

### 3) IMDb : Movie Data, Actor Data

- Predict how many people will like a movie.
- What is gross?
- Social network of actors - (Six degrees of separation)

### 4) NYC Taxi Cab data

- Freedom of Information Act Request (FOIA) 으로부터 구할 수 있는 data.
- How far do they travel?

### 5) Google Ngrams

- Presents an annual time series of the frequency of every "popular" word / phrase with 1 to 5 words occurs in scanned books.  
↑ appear more than 40 times
- What is the lifespan of fame and technologies?  
Is it increasing or decreasing?

### 6) Google Trends

- search term

- example ?

# Properties of Data

## 1) Structured vs. Unstructured

Tweets from Twitter  
Movie reviews (free text)  
X-ray image  
YouTube videos

} unstructured data

Sales record spreadsheet  
Financial Time series

} structured data

## 2) Data Types

### - Nominal (Categorical) (N)

: labels.  $= \neq$

No order. 같다 / 다르다만 구분 가능. Index.

### - Ordinal (O)

: ordered.  $= \neq > <$

### - Quantitative (Q)

#### - Q-Interval

: location of zero arbitrary.  $= \neq > < + -$

only Intervals (distance) can be compared. Dates, Locations..

#### - Q-Ratio

: zero fixed. origin is important.  $= \neq > < + - \times \div$

can measure ratios & proportions. Length. Mass. Temperature..

## • Classification vs. Regression

### - Classification

: Assign a label to an item from a discrete set of probabilities.

### - Regression

: Forecast a numerical quantity (amount). the value is continuous.



## Lecture 2. Mathematical Preliminaries

### 2.1 Probability

- Probability theory provides a formal framework for reasoning about the likelihood of events.
- experiment: a procedure which yields one ~~se~~ of a set of possible outcomes. rolling <sup>two</sup> dice
- sample space S: the set of possible outcomes of an experiment.  $S = \{(1,1), (1,2), \dots, (6,6)\}$
- event E: a specified subset of the outcomes of an experiment. sum of dice = 7  
 $E = \{(1,6), (2,5) \dots (6,1)\}$
- probability of an outcome s:  $p(s)$  where  $\begin{cases} 0 \leq p(s) \leq 1 \\ \sum_{s \in S} p(s) = 1 \end{cases}$
- probability of an event E:  $p(E) = \sum_{s \in E} p(s) = 1 - p(\bar{E})$
- random variable V: numerical function on the outcomes of a probability space.  
확률공간 S를 정의역으로 하고, 실수를 공역으로 가지는 함수로 정의.  
확률변수를 사용하면 모든 확률을 실수에 대응시켜 표현할 수 있다.  
따라서 확률공간 S를 수리적인 위에 표현할 수 있고 사건 E는 이 수리적인 상 일부 구간으로 표현된다.  
 $V(a,b) = a+b$   
random variable V를 sum of values of two dice로 정의했을 때이다.
- expected value: mathematical expected value of a random variable is defined as  
$$E(V) = \sum_{s \in S} p(s) \cdot V(s)$$

#### 2.1.1 Probability vs. Statistics

- probability: predicting the likelihood of future event
- statistic: analyzes the frequency of past event  
applied to mathematics trying to make sense of observations in the real world.  
↳ past events

#### 2.1.2 Compound Events and Independence

- probability of a compound event is combination of probability of two or more simple events.
- The events A and B are independent if and only if  
$$P(A \cap B) = P(A) \times P(B)$$
- Independence (zero correlation): good to simplify calculations but  
bad for prediction (!!)

### 2.1.3 Conditional Probability

$$P(A|B) = \frac{P(A \cap B)}{P(B)} \quad \text{B가 일어났을 때 A가 일어날 확률}$$

- conditional probabilities get interesting only when events are NOT independent.

$\therefore A, B$ 가 independent 한 경우,  $P(A|B) = \frac{P(A \cap B)}{P(B)} = \frac{P(A)P(B)}{P(B)} = P(A) \leftarrow$  의미 없다!

#### - Bayes Theorem

: Important tool which reverses the direction of the dependencies.

$$P(B|A) = \frac{P(A \cap B)}{P(A)} = \frac{\frac{P(A \cap B)}{P(B)} \cdot P(B)}{P(A)} = \frac{P(A|B) P(B)}{P(A)}$$

$\Rightarrow \boxed{P(B|A) = \frac{P(A|B) P(B)}{P(A)}}$  구하기 쉬운 것을 이용해서 구할 수 있다!

구하기 어려울 때

ex)  $P(\text{spam} | \omega) = \frac{P(\omega | \text{spam}) \cdot P(\text{spam})}{P(\omega)}$

word 보고 spam 판단  $\rightarrow$  이진.

spam 보고 word 찾기  $\rightarrow$  사용.

### 2.1.4 Probability Distributions

- Random variables are numerical functions where the values are associated with probabilities of occurrence.

- RVs can be represented by their probability density function (pdf).

- pdf:  $\begin{cases} x\text{-axis} : \text{range of values the random variable can take on} \\ y\text{-axis} : \text{probability of that given value.} \end{cases}$

$$P(k=X) = \frac{h(k=X)}{\sum_x h(x=X)}$$

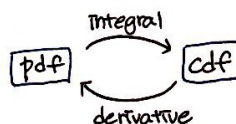
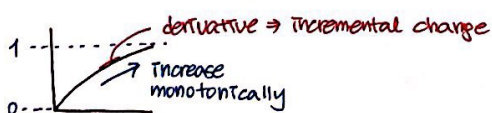
- pdf 는 histogram 과 유사하지만 y 축이 frequency 가 아니라 probability 라는 차이가 있다.

- RVs can also be represented by their cumulative density function (cdf)

$\rightarrow$  may present misleading view of growth rate. (Apple case)

- cdf: running sum of the probabilities in the pdf. it reflects  $\underbrace{P(X \leq k)}_{\text{cdf}}$  instead of  $\underbrace{P(X = k)}_{\text{pdf}}$

$$C(X \leq k) = \sum_{x \leq k} P(X=x)$$



: exact same information

① representing data

② reducing data

2.2 Descriptive Statistics → data에 대한 summary 제공, 대표값을 잘 제공.

- Descriptive statistics provides way to capture the properties of a given data set / sample

1) Central tendency measure : describe the center around which data is distributed

2) Variation or variability measures : describe data spread, i.e. how far the measurements are from the center.

### 2.2.1 Centrality Measures

#### 1) Arithmetic Mean

$$\mu_x = \frac{1}{n} \sum_{i=1}^n x_i$$

: Mean is meaningful for symmetric distributions without outliers, like height or weight.

( $\mu$ 보다 작은 item  $\downarrow$ )

$\approx$  ( $\mu$ 보다 큰 item  $\downarrow$ )

#### 2) Median : middle value.

(중간값)

Median is better for skewed distributions or data with outlier, like wealth.

#### 3) Mode : most frequent element in the data set.

(가장 많이 나타나는 값)

Mode is often ~~is~~ NOT close to the center.

text data에서 주로 사용. 유사한 특징으로 grouping 할 때 mode는 해당 cluster를 대표하기 좋음.

#### 4) Geometric Mean

$$\left( \prod_{i=1}^n a_i \right)^{\frac{1}{n}} = \sqrt[n]{a_1 a_2 \dots a_n}$$

nth root of product of n values

: geometric mean  $\leq$  arithmetic mean  
always

: geometric mean is sensitive to zero. (0이 아닌 0에 가까운 값의 영향이 큼!)

: Geometric means make sense with ratios.  $\frac{1}{2}, 2$ 의 geometric mean = 1

### - Aggregation as Data Reduction

Representing a group of elements by a new derived element, like mean, min, count,

sum reduces a large amount of dataset to a small summary statistic.

Reduce by  
a new derived element

can become features when taken over natural groups or clusters in the full data set.



### 2.2.2 Variability Measures

#### - Standard Deviation

$$\sigma = \sqrt{\frac{\sum_i^n (x_i - \bar{x})^2}{n-1}}$$

← population size  $n$ , sample size  $n-1$ .

for large  $n$ ,  $n \sim (n-1) \therefore$  doesn't matter.

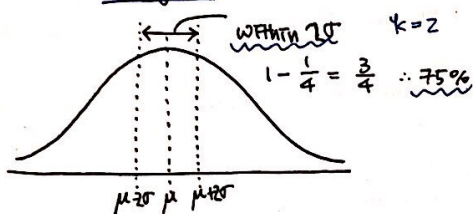
$$\text{Variance} = \sigma^2$$

#### - mean, standard deviation

together  $\Rightarrow$  characterize any distribution well

#### - parameterizing distributions

"Regardless of how data is distributed", at least  $(1 - \frac{1}{k^2})$ th of the points must lie within  $k$  sigma of the mean.



$\Rightarrow$  tighter bounds apply for normal distribution.  
 normal size  $\pm 2\sigma$  or 95%

### 2.2.3 Interpreting Variance

- Repeated observations of the same phenomenon do not always produce the same results, due to random noise or error.

- sampling error: when observations capture unrepresentative circumstances  
 e.g. measuring rush hour traffic as on weekend as well as week days.

- measurement error: limits of precision inherent in any sensing device.

- signal to noise ratio: the degree to which a series of observations reflects a quantity of interest as opposed to data variance.

- hard to measure 'signal to noise ratio', because much of what you see is often variance due to sampling error and measurement error.

1) stock market: measuring the relative "skill" of different stock market investors.

2) sports performance: 선수의 실력이 시즌에 따라 변하는 것.  $\rightarrow$  genuine difference..?

- 3) Model performance : simple model ~ complex model  
 ↳ larger variance.  
 ↳ more generalizable

## 2.3 Correlation Analysis

- Two factors are correlated when values of  $x$  has some predictive power on the value of  $y$ .
- correlation coefficient  $r(X, Y)$  : statistic that measures the degree to which  $Y$  is a function of  $X$ , and vice versa.
- correlation ranges from -1 to 1.
  - 1 means anti-correlated
  - 0 means no relation
  - 1 means fully correlated

### 2.3.1 Correlation Coefficient : Pearson and Spearman Rank

There are two primary statistics used to measure correlation. These different statistics are appropriate in different situations.

#### 1) The Pearson Correlation Coefficient

← more prominent

$$r = \frac{\sum_{i=1}^n (X_i - \bar{X})(Y_i - \bar{Y})}{\sqrt{\sum_{i=1}^n (X_i - \bar{X})^2} \sqrt{\sum_{i=1}^n (Y_i - \bar{Y})^2}} = \frac{\text{Cov}(X, Y)}{\sigma(X)\sigma(Y)}$$

defines the degree to which a linear predictor of the form  $y = m \cdot x + b$  can fit the observed data.

⇒ linear predictor 71 011 739 911 111.

e.g.  $y = 1x + 1$  ← correlation 0.48.

-  $X, Y$  are correlated (strongly)

$\left\{ \begin{array}{l} X, Y \text{ 739 739} \Rightarrow \text{positive} \\ \text{739} \Rightarrow \text{negative} \end{array} \right.$

$(X_i - \bar{X} > 0 \text{ 739 } Y_i - \bar{Y} > 0)$

-  $X, Y$  are uncorrelated

positive & negative terms should occur with equal frequency, offsetting each other.

⇒ value becomes zero.

numerator 739. (determining sign of correlation) → named as covariance.

$$\text{Cov}(X, Y) = \sum_{i=1}^n (X_i - \bar{X})(Y_i - \bar{Y}) / n$$

denominator : reflects the amount of variance in the two variables, as measured by their standard deviation.

covariance between  $X$  and  $Y$  potentially increases with the variance of these variables

→  $r \approx 1 \sim 1 \approx 739$ .



## 2) Spearman Rank Correlation Coefficient

better with non-linear relationships and outliers.  
less sensitive to outlier elements than Pearson.

- spearman rank correlation coefficient : counts the number of pairs of input points which are out of order.

$$\rho = 1 - \frac{6 \sum d_i^2}{n(n^2 - 1)}$$

only cares about rank not the actual value

counts the number of disordered pairs, not how well the data fits.  $-1 \leq \rho \leq 1$

$$d_i = \text{rank}(x_i) - \text{rank}(y_i)$$

$\text{rank}(x_i)$  : the rank position of  $x_i$  in sorted order among all  $x_i$ .

$$1 \leq \text{rank}(x_i) \leq n$$

$p(x_1, y_{\max}) \rightarrow p(x_1, \infty)$  : spearman 은 평균값  $x$  ( $\because \text{rank}(y)$  변하지 않아서!)  
pearson  $\Rightarrow$  so crazy!

### 2.3.2 The Power and Significance of Correlation

- Correlation becomes more impressive the more points it is based on.

#### 1) Strength of correlation : $R^2$

square of correlation coefficient

$r^2$  estimates the fraction of the variance in  $Y$  explained by  $X$  in a simple linear regression.

$\rightarrow Y$  분산의  $r^2$  만큼만  $X$ 가 설명할 수 있다.

linear fit  
Let  $f(x) = mx + c$  be predicted value of  $y$  from  $x$ .

Then, the residual value  $r_i = y_i - f(x_i)$ . (mean of  $r = 0$ )

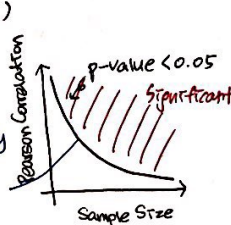
If  $x, y$  are perfectly correlated, there will be no residual error.  $\therefore V(r) = 0$

If  $x, y$  are totally uncorrelated, the fit should contribute nothing  $\therefore V(y) \approx V(r)$

for good linear fit  $f(x)$ ,  $V(r) < V(y)$

$$\Rightarrow 1 - r^2 = V(r) / V(y)$$

the level of correlation necessary to be statistically significant



Even small correlations become significant ( $p\text{-value} < 0.05$ ) with large enough sample sizes.

#### 2) Statistical Significance

The statistical significance of a correlation depends upon its "sample size  $n$ " as well as  $r$ .

Traditionally, a correlation of  $n$  points is significant (if) there is an  $\alpha \leq \frac{1}{20} = 0.05$

chance that we would observe a correlation as strong as  $r$  in any random set of  $n$  points. \* large # of weak but independent correlations may together have strong predictive power.

random 요소 (무연관)  $r$  만큼의 correlation을 관찰하는 확률이 0.05 미만일 때  $n$  개의 point 에러 관찰한 correlation  $r$  이 significant 하다.  $\hookrightarrow p\text{-value}$

### 2.3.3 Correlation Does Not Imply Causation

- correlation does not mean causation. e.g. police, crime

### 2.3.4 Detecting Periodicities by Autocorrelation

- Autocorrelation: comparing a sequence to itself. → important concept in predicting future events.
- Autocorrelation function: series of correlations for all  $1 \leq k \leq n-1$
- Time series data often exhibits cycles which affect its interpretation.  
A "cycle of length  $k$ " can be identified by unexpectedly large autocorrelation between  $S[t]$  and  $S[t+k]$  for all  $0 < t < n-k$ .
- Computing the lag- $k$  autocorrelation takes  $O(n)$ , but the full set can be computed in  $O(n \log n)$  via the Fast Fourier Transform (FFT).

## 2.4 Logarithms

- logarithm: inverse exponential function.

$$y = \log_b x \Rightarrow b^y = x$$

- ① Summing logs of probabilities is more numerically stable than multiplying them.

$$\prod_{i=1}^n p_i = b^p \quad \text{where } p = \sum_{i=1}^n \log_b(p_i)$$

ratio, power law 에 대해

log 취하면 → good ~~distribution~~ analysts

### ② Logarithms and Ratios ✓

- 비율이 2배 증가하거나, 1/2배로 감소하는 상황. ratio 는 4배 차이난다.

log 를 취하면 +1.0, -1.0 으로 같은 크기의 차이로 표현할 수 있다.

- ratio 만 보면 outlier처럼 보이는 것 → log 취하면 보면 outlier 가 아닐 수 있다!

### ③ Logarithms and Power Laws ✓

- power law distribution 에 log 취하면 traditional distribution 으로 보일 수 있다.

### ④ Normalizing Skewed Distributions

